

Why Sequence Models?

- Examples of Sequence data :-

- ① Speech recognition :-  → "The quick brown fox jumps over the lazy dog".
- ② Music generation :- ϕ → 
- ③ Sentiment classification :- "bad movie" → 1★
- ④ DNA Sequence Analysis → AGCCCTTACCG
- ⑤ Machine translation → "Hermosa" → "Beautified".
- ⑥ Video activity recognition →  → "Running".
- ⑦ Name entity recognition → "Yesterday Harry Potter met Hermoine granger."

NOTATIONS

x : "Harry Potter & Hermoine granger invented a new spell".

Name entity recognition

y : 1 1 0 1 1 0 0 0 0

- input is sequence of 9 words, 9 sets of features to represent these 9 words.

$x^{(t)}$, $x^{(1)}$, $x^{(2)}$ --- $x^{(9)}$ $T_x = 9$

"temporal Sequence"

$y^{(t)}$, $y^{(1)}$, $y^{(2)}$ --- $y^{(9)}$ $T_y = 9$

$x^{(i)(t)}$, $T_x^{(i)}$

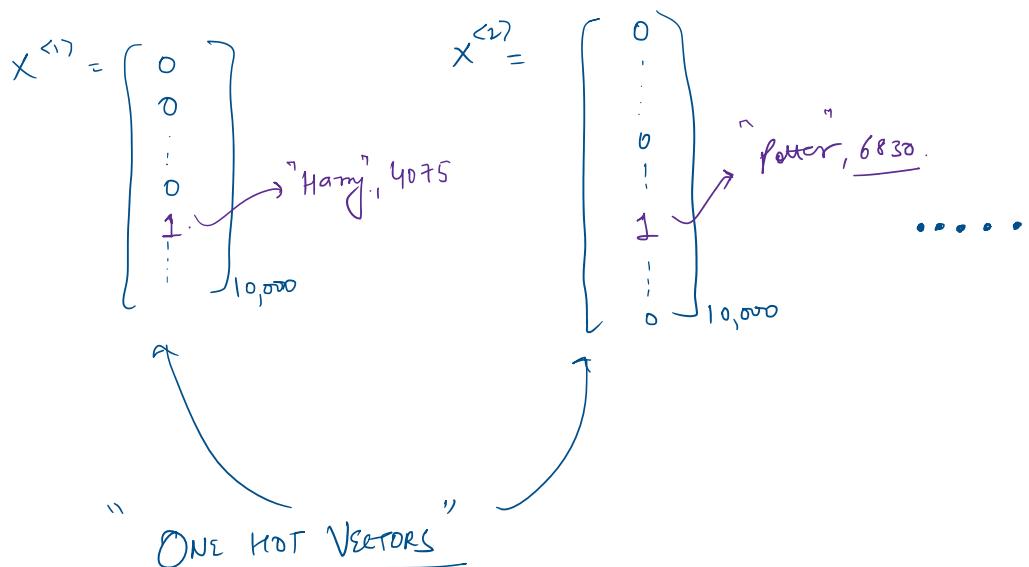
$$x^{(i)<\tau}, \quad T x^{(i)}$$

$$y^{(i)<\tau}, \quad T y^{(i)}.$$

Vocabulary	dictionary
a	1
and	2
Harry	367
Potter	9075
Zebra	6830
	10,000

(Assuming).

One Hot representations



Recurrent Neural N/w MODEL.

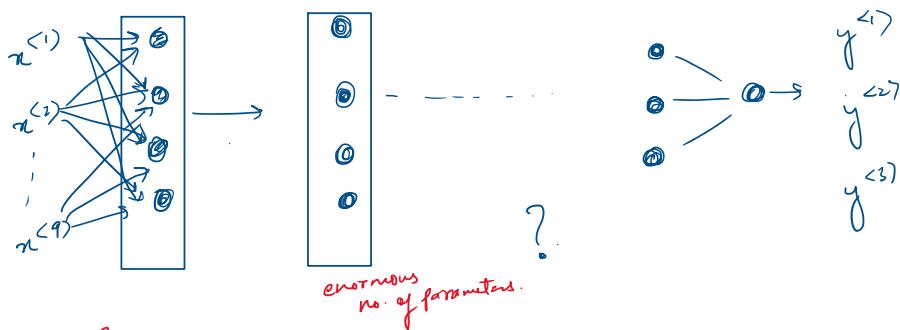
Why not a Standard Neural N/w ?

$$x^{(i)} \rightarrow \boxed{\oplus}$$

$$\boxed{\oplus}$$

$$y^{(i)}$$

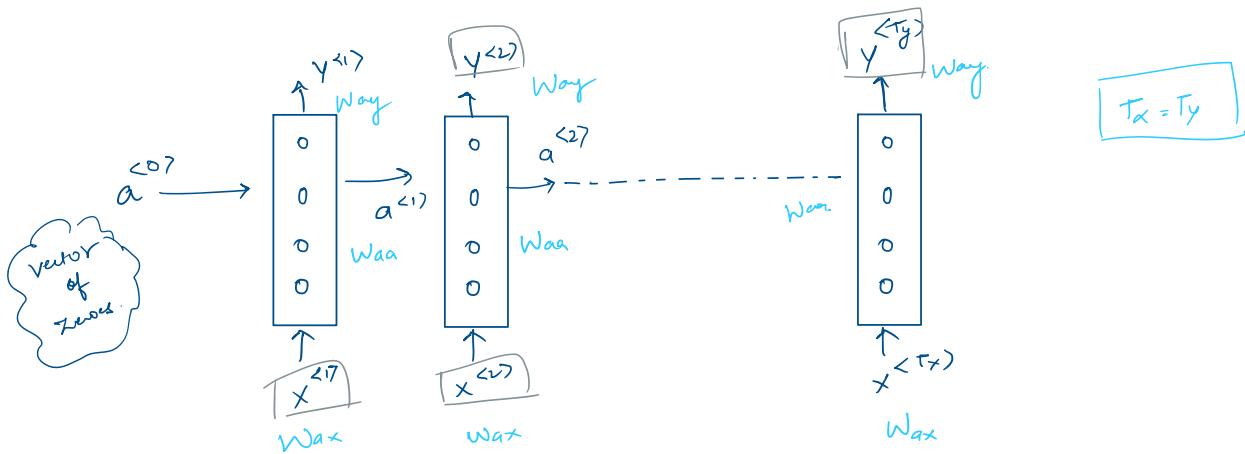
0



Problems!

- ① I/P and O/P can be different layers in diff example.
- ② doesn't share features learned across different portion of text.

Recurrent Neural NW

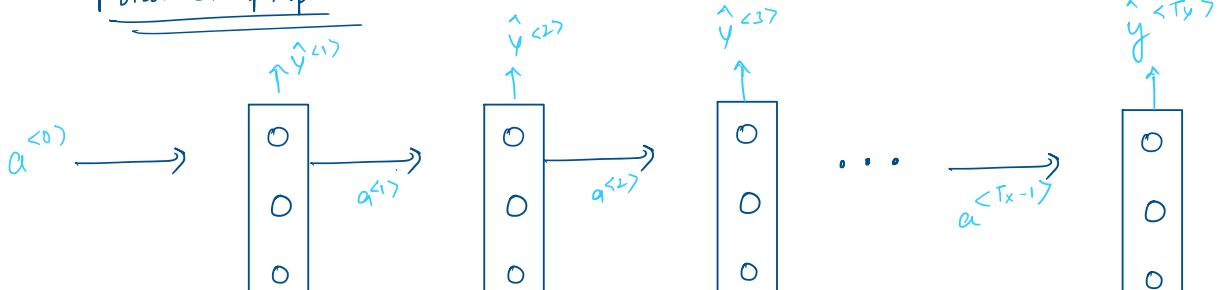


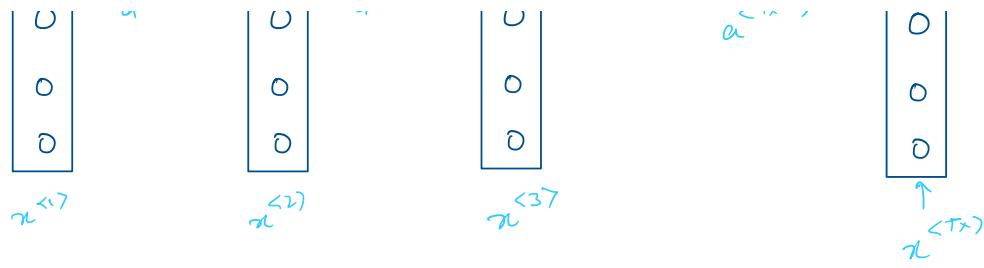
> While making predictions $y^{(t)}$, it also uses $x^{(t)}$ and $x^{(t-1)}$
 or
 (earlier in the sequence.)

ex: He said "Teddy Roosevelt was a great president." (✓)

He said "Teddy Bears are on Sale". (✗)

Forward Prop.





$$a^{<0>} = \vec{0}.$$

$$a^{<1>} = g_1(W_{aa} \cdot a^{<0>} + W_{ax} \cdot x^{<1>} + b_a)$$

$$\hat{y}^{<1>} = g_2(W_{ya} \cdot a^{<1>} + b_y)$$

$\rightarrow \text{tanh}(\cdot) / \text{ReLU}$

$\rightarrow \text{Sigmoid}$.

$$a^{<t>} = g(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya} \cdot a^{<t>} + b_y).$$

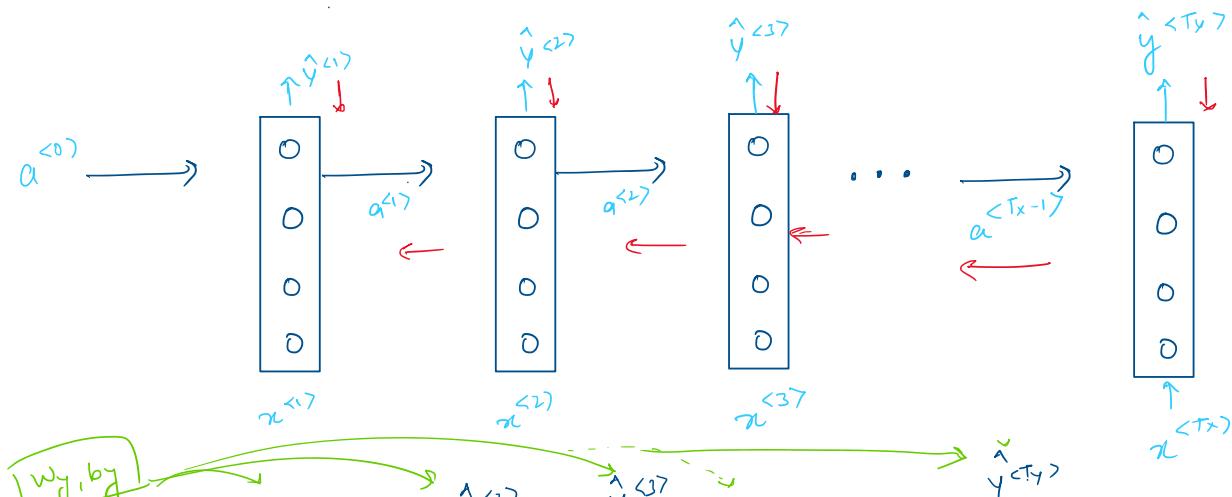
Simplified RNN Notation

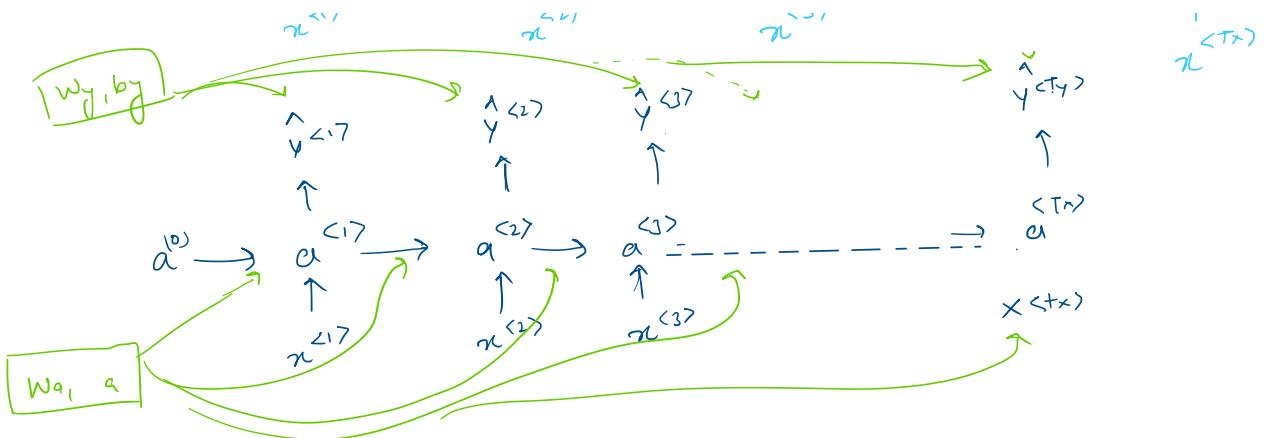
$$a^{<t>} = g(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a) \iff g(W_a(a^{<t-1>} + x^{<t>})) + b_a$$

$$\hat{y}^{<t>} = g(W_{ya} \cdot a^{<t>} + b_y).$$

$$\begin{bmatrix} W_{aa} & | & W_{ax} \end{bmatrix} \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = \boxed{W_a(a^{<t-1>} + x^{<t>})}$$

BACKPROP. THROUGH TIME



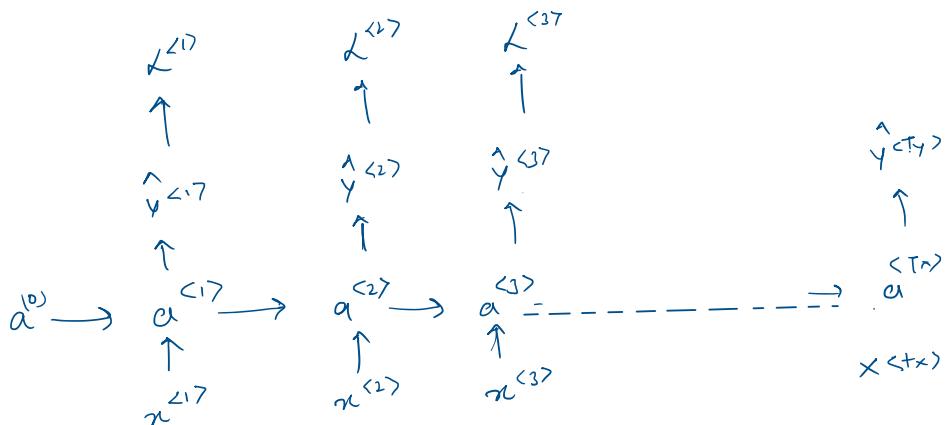


Loss fx.

$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1-y^{<t>}) \log (1-\hat{y}^{<t>})$$

"Cross entropy loss"

$$L(y, \hat{y}) = \sum_{t=1}^{Tx} L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$



Different Types Of RNN.

[when $T_x = T_y$]

Type of RNN	Illustration	Example
One-to-one $T_x = T_y = 1$		Traditional neural network

One-to-many $T_x = 1, T_y > 0$		Music generation
Many-to-one $T_x > 1, T_y = 1$		Sentiment classification
Many-to-many $T_x = T_y$		Name entity recognition
Many-to-many $T_x \neq T_y$	 From https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks#overview	Machine translation

Stanford University CS230

"Encoder" "Decoder"

LANGUAGE MODEL AND SEQUENCE GENERATION.

Speech Recognition

"The apple and pear salad."

pear / pair. ?

- $P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$ ✓

$$P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$$

$$P(\text{pair}) = ?$$

$$y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(T_y)}$$

Language Building with RNN.

$T_{in} \dots \text{[yellow]} \dots T_{out}$

v v v

Training set \Rightarrow Large **Corpus** of english text.

Collection

Example Sentence :-

Cats average 15 hours of sleep a day.

① Tokenize (Create vocabulary/dictionary).

$y^{<1>} \quad y^{<2>} \quad y^{<3>} \dots \dots \dots \quad y^{<T_y>} \quad <\text{EOS}>$

append an extra token at the end of sequence $<\text{EOS}>$
end of sentence.

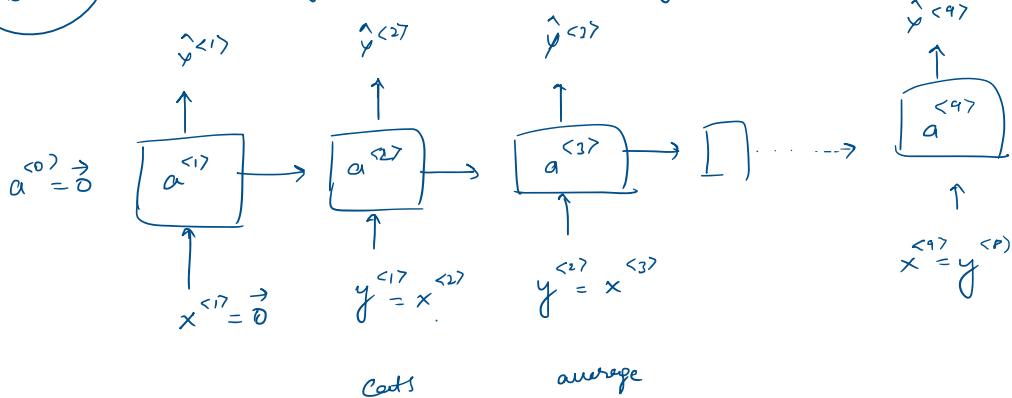
"The Egyptian Mau is a breed of cat."

if the word is not in vocabulary
replace with $\langle \text{UNK} \rangle$ token.
 \nearrow unknown word.

RNN Model

choose example

• Cats average 15 hours of sleep a day $<\text{EOS}>$



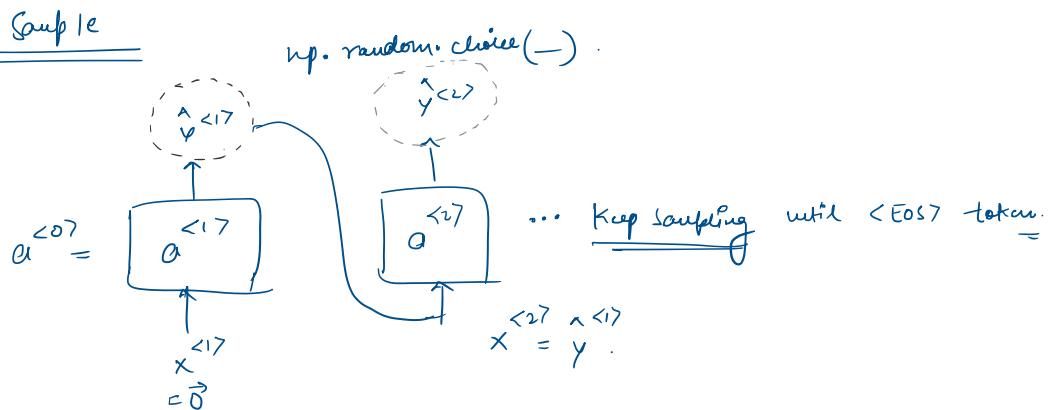
$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

$$\Rightarrow L = \mathcal{L}(\hat{y}^{<t>}, y^{<t>})$$

SAMPLING NOVEL SEQUENCES.

$$P(y^{<1>} | y^{<2>} | \dots | y^{<T>})$$

Sample

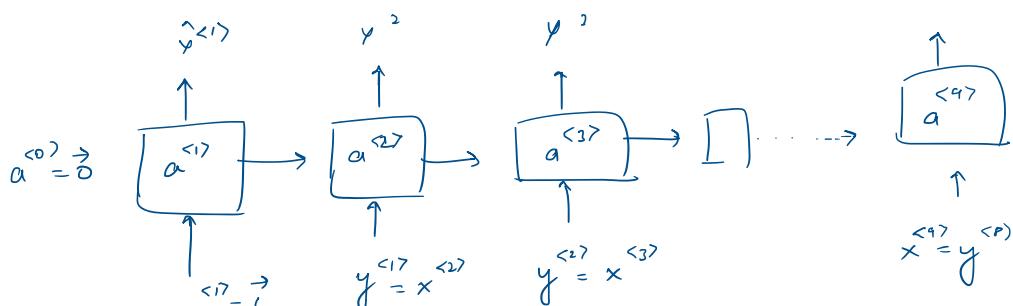


Character level words -

Vocabulary = $[a, b, c, \dots, z, A-Z, ., , :, ;, " - \dots]$

Very expensive.

VANISHING GRADIENTS WITH RNN'S



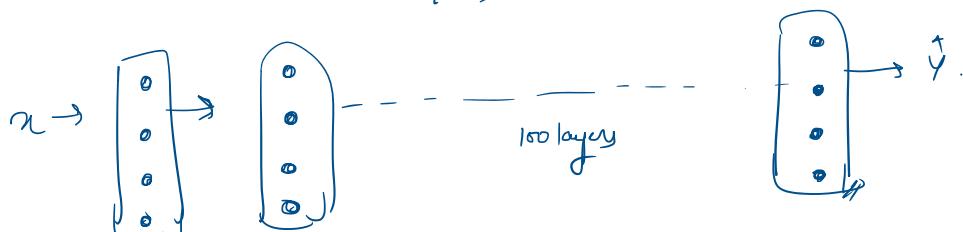
sentence : "The cat which already ate -----, was full."

Cat is singular so

Cat : was

Cats : were.

ex. in deep neural net \Rightarrow
(100 layers)





forward prop. $\xrightarrow{\quad}$ Backprop.

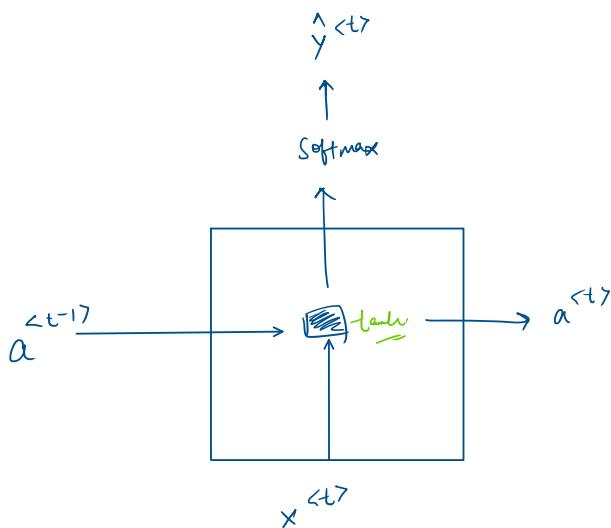
- ④ if the layer is very deep, we will have hard time calculating gradient, affecting weights of earlier layers.

Vanishing / Exploding gradients $\xrightarrow{\quad}$ by tuning weight.

\rightarrow gradient clipping.

GATED RECURRENT UNIT. (GRU)

$$a^{<t>} = \tanh(g(W_a [a^{<t-1>}, x^{<t>}] + b_a))$$



"The cat which already ate . . . , was full."
 {cat : was } $\xrightarrow{\quad}$ {cats : were }

C = memory cell. $\left. \right\} \text{Provide memory to remember that the "cat" was singular/plural}$

c = memory cell. $\left\{ \begin{array}{l} \text{Provide memory to remember that the "cat" was singular/plural} \\ \text{at time } t \quad \text{for further consideration.} \end{array} \right.$

$$c^{<t>} \xrightarrow{o/p} a^{<t>}$$

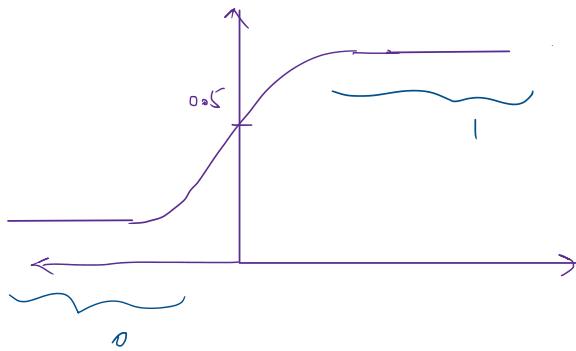
$$\tilde{c}^{<t>} = \tanh(w_c [c^{<t-1>} \times^{<t>}] + b_c).$$

Candidate for replacing $c^{<t>}$.

update gate

$$\gamma_u = (0, 1) = \sigma(w_u [c^{<t-1>} \times^{<t>}] + b_u).$$

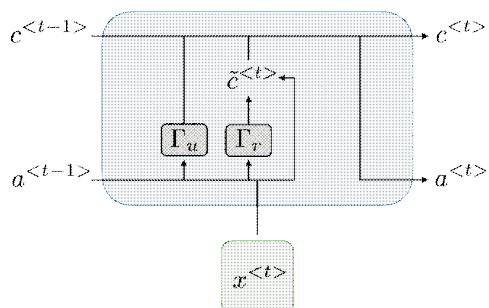
(gamma)



$\left\{ \begin{array}{l} \text{singular : 1} \\ \text{plural : 0} \end{array} \right.$

$$c^{<t>} = \gamma_u * \tilde{c}^{<t>} + (1 - \gamma_u) * c^{<t-1>}$$

FULL GRU



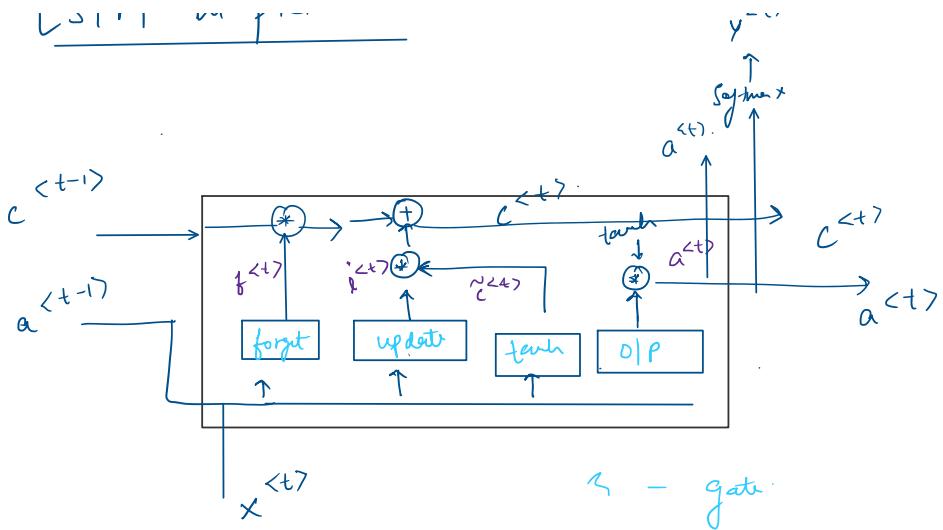
$$\begin{aligned}
 h & C^{t+1} = Y_u * \tilde{C}^{t+1} + (1 - Y_u) * C^{t+1} \\
 \tilde{h} & \tilde{C}^{t+1} = \tanh \left(w_c [Y_r * C^{t+1}, x^{t+1}] + b_c \right) \\
 r & Y_r = \sigma(w_r [C^{t+1}, x^{t+1}] + b_r) \\
 u & Y_u = (0, 1) = \sigma(w_u [C^{t+1}, x^{t+1}] + b_u) \\
 a^{t+1} & = C^{t+1}
 \end{aligned}$$

LONG-SHORT TERM MEMORY. (LSTM)

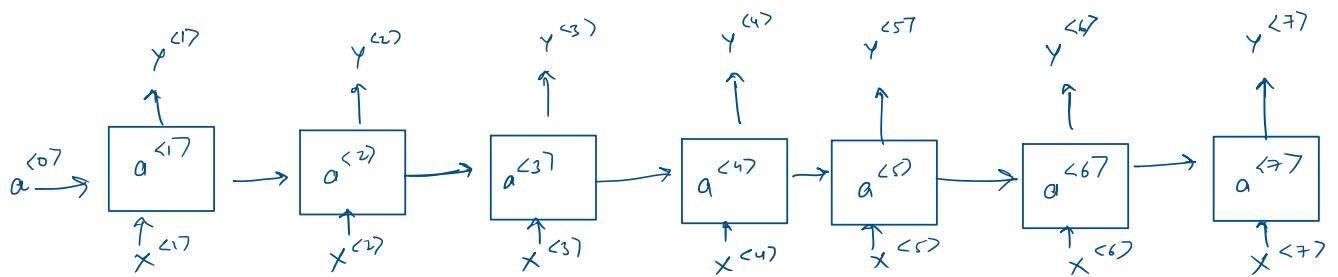
$$\begin{aligned}
 C^{t+1} &= Y_u * \tilde{C}^{t+1} + (1 - Y_u) * C^{t+1} \\
 \tilde{C}^{t+1} &= \tanh \left(w_c [a^{t+1}, x^{t+1}] + b_c \right) \\
 Y_r &= \sigma(w_r [C^{t+1}, x^{t+1}] + b_r) \\
 Y_u &= (0, 1) = \sigma(w_u [a^{t+1}, x^{t+1}] + b_u) \\
 Y_f &= \sigma(w_f [a^{t+1}, x^{t+1}] + b_f) \\
 &\downarrow \\
 & \text{"forget gate"} \\
 V_o &= \sigma(w_o [a^{t+1}, x^{t+1}] + b_o) \\
 &\uparrow \\
 & \text{"Output"} \\
 C^{t+1} &= [Y_u * \tilde{C}^{t+1}] + [Y_f * C^{t+1}] \\
 a^{t+1} &= Y_o * \tanh(C^{t+1})
 \end{aligned}$$

LSTM in picture

y^{t+1}
 \uparrow
 softmax



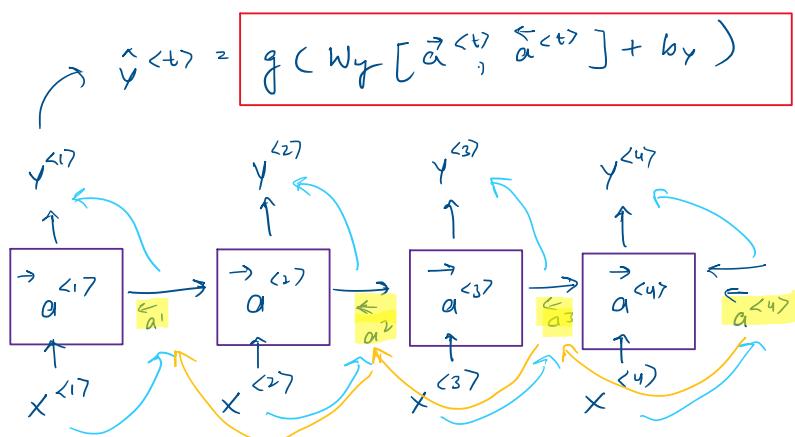
BIDIRECTIONAL RNN. (BRNN)



He said " Teddy Bears are on sale"

He said "Teddy Roosevelt was a great President"

- Unidirectional RNN.

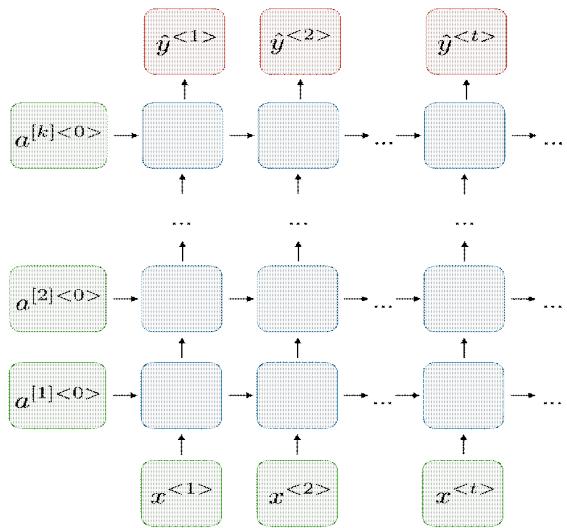


A word sentence.

Acyclic graph

" BRNN with LSTM "

DEEP RNN



MI INTRODUCTION To WORD-EMBEDDINGS.

WORD REPRESENTATION.

$$V = [q, \text{aeron}, \dots, \dots, \text{Zulu}, \text{UNK}]$$

1 hot representation

Man (5391)	Women (9853)	King (4914)	Queen (7157)
$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ -1 \\ 0 \\ \vdots \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{pmatrix}$	- - - -	- - - -

Featured Representation

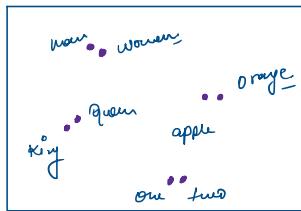
	Man	Woman	King	Queen	Apple	Orange
gender	-1	1	-0.95	0.97	0.50	0.01
royalty	-0.01	-0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.61	-0.03	-0.2
foot						
size						
:						
---			-	-	-	-
e ₍₃₉₎			-	-	-	-

" I want a glass of Orange juice.

" I want a glass of apple juice.

* Orange and apple are highly correlated ("friend").

Visualizing Word Embeddings

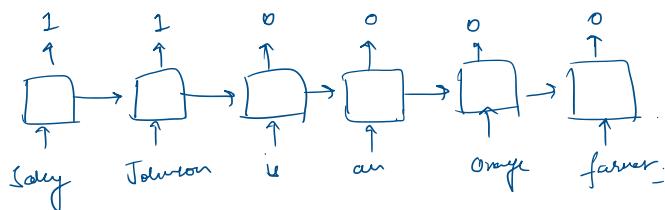


t-SNE.

USING WORD EMBEDDINGS

" Name Entity Recognition "

Sally Johnson is an Orange farmer.



" Robert is an durian farmer "
 Robert = a
 durian = Cultivator.
 farmer = Cultivator.

"
 ⇒ Cultivator & farmer are highly correlated."
 "

Transfer Learning

- * > Learn word embedding from large text-corpus.
- * > Transfer embeddings to new task with smaller training set.
- > Continue to fine tune the word embeddings with new data.

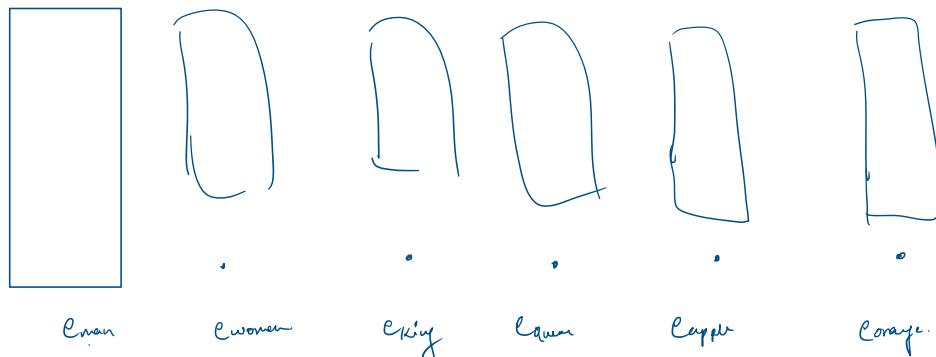
PROPERTIES OF WORD-EMBEDDINGS

Analogies.

Man : Women :: King : ?

Man	Women	King	Queen	Apple	Orange
(5391)	(9853)	(4914)	(7153)	(456)	(6257)

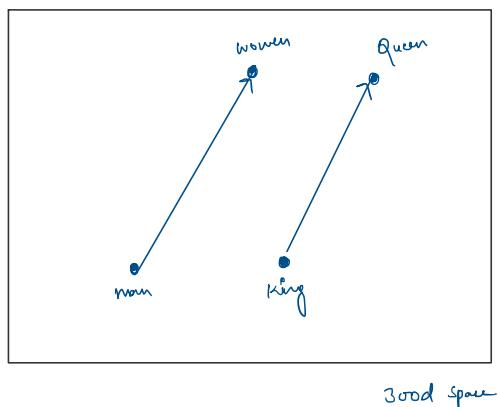
gender
royalty
Age
Food



$$e_{\text{man}} - e_{\text{woman}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad e_{\text{king}} - e_{\text{queen}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - [\dots]$$

$$\Rightarrow e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}}$$

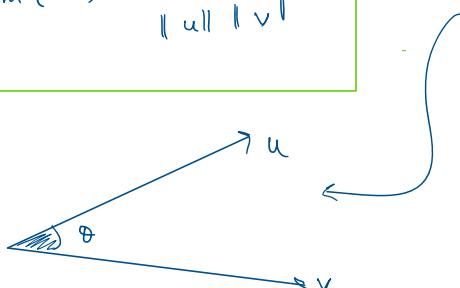


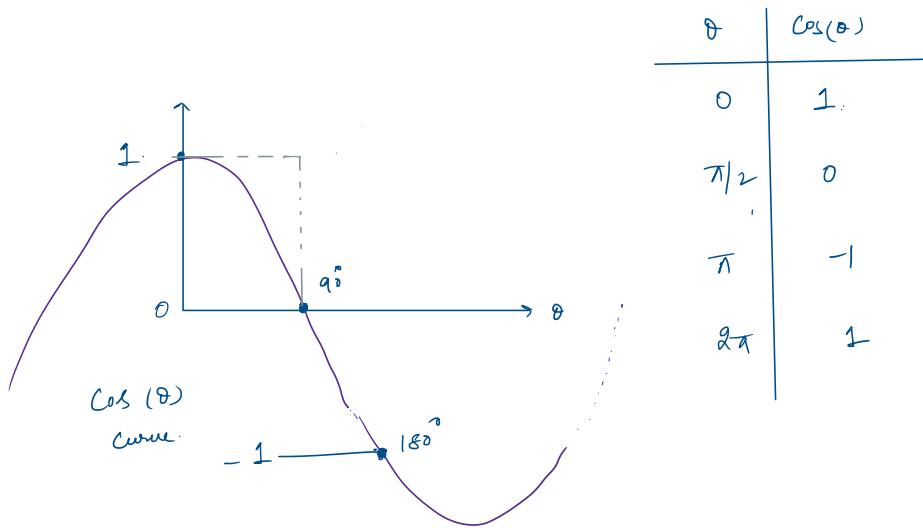
find word w , s.t:

$$\max \text{Sim.}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$

- find a similarity fx.
 \Rightarrow Cosine Similarity.

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|} = \cos(\theta)$$





Other similarity fx:

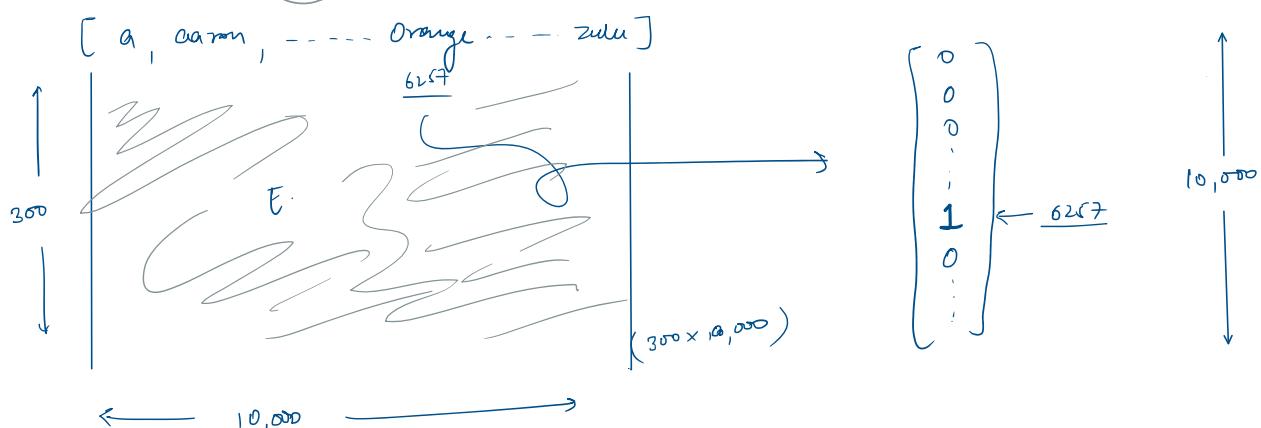
- Euclidean distance. $\| u - v \|^2$

India : delhi :: Kenya : "Nairobi"

EMBEDDING MATRIX.

embedding matrix (E)

($300 \times 10,000$) \rightarrow if you have 10,000 word dictionary



$$E \cdot O_{6257} = \begin{bmatrix} \text{apple} \\ \text{orange} \\ \text{carrot} \end{bmatrix}_{300, 1} = e_{6257}$$

one-hot vector.

$$\boxed{E \cdot O_j = e_j}$$

Embedding for word j .

Embedding for word.

* It is not efficient to use multiplication of matrix, as the cost is very high.

MD LEARNING WORD EMBEDDINGS : Word2Vec and Glove

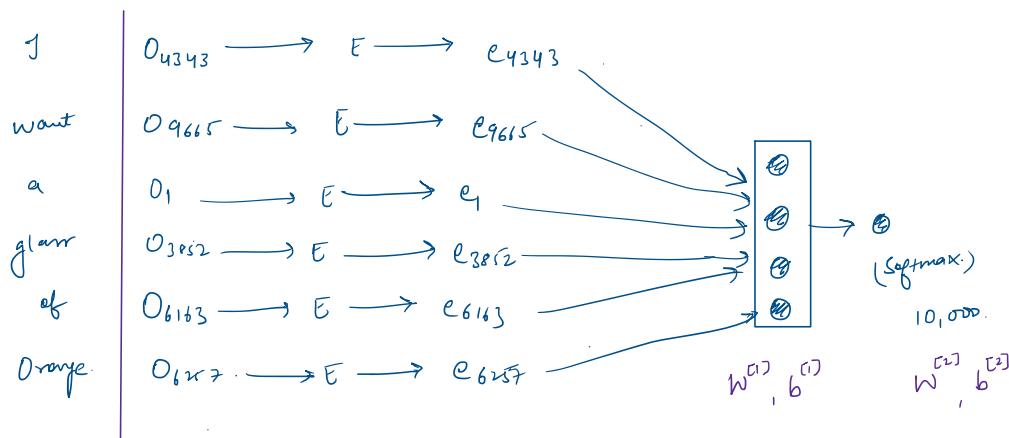
LEARNING WORD EMBEDDING

Natural Language Model.

"I want a glass of Orange _____."

(4343) (9665) (3852) (6163) (6257)

.....
index in vocab.



WORD2VEC

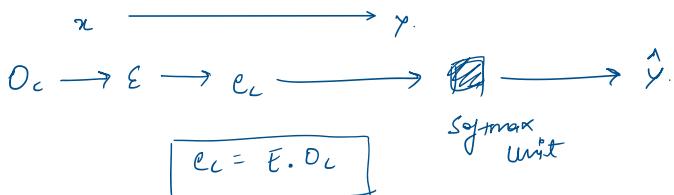
Skip-grams : I want a glass of Orange juice along with my cereal.

Context	Target
randomly	juice
selected	glass
	my

Model

$$\text{vocab-size} = 10,000 \text{K}$$

Context c ("Orange") \longrightarrow Target t ("juice")
 (6257) (4834)



Softmax : $p(t|c) = \frac{e^{\theta_t^T \cdot e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T \cdot e_c}}$

θ_t = parameter associated with O/P t
 (class of t , being a label).

$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{1000} [y_i \cdot \log \hat{y}_i]$$

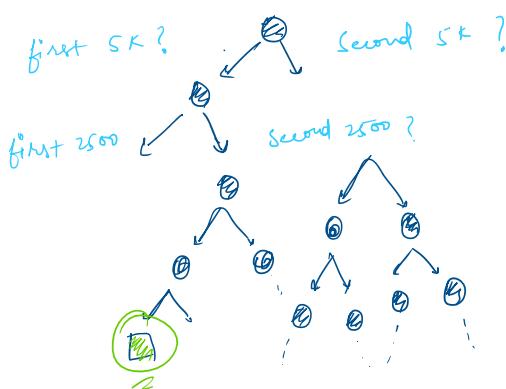
Problems with Softmax Classification.

- Computational Speed.

$$p(t|c) = \frac{e^{\theta_t^T \cdot e_c}}{\sum_{i=1}^{10,000} e^{\theta_i^T \cdot e_c}}$$

every time it will calculate 10k times,
 at higher values, it will compute slower.

\Rightarrow Hierarchical softmax.



NEGATIVE SAMPLING.

Predict (Context, target) pair

Context	target	<u>Probability</u>
Orange	juice	≈ 1
Orange	king	≈ 0
Orange	book	≈ 0

K

find p with random K target and fill them (0).

> Create a supervised learning algorithm, given (Context, target) predict probability.

$$K = \begin{cases} 5-20, & \text{for small datasets} \\ 2-5, & \text{for larger datasets} \end{cases}$$

Softmax Model \rightarrow

$$p(t|c) = \frac{e^{\theta_t^T \cdot e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T \cdot e_c}}$$

Context	word	target
- - -	- - -	1
- - -	- - -	0
- - -	- - -	0
- - -	- -	0

use Logistic Regression Model.

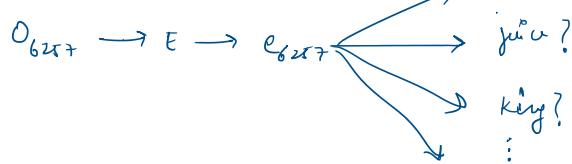
$$p(y=1 | c, t) = \sigma(\theta_t^T \cdot e_c)$$

Sigmoid

"Orange"

6257

e_{6257}



1, real ex.
k, random ex.

Selecting -ve examples?

$$P(w_i) = \frac{f(w_i)^{1/4}}{\sum_{j=1}^{10,000} f(w_j)^{1/4}} \approx \frac{1}{|V|}$$

uniform distribution raised to $1/4^{th}$.

$$\sum_{i=1}^{10,000} f(w_j)^{3/4}$$

|V|

to $S^T T^{-1}$

GLOVE WORD VECTORS

global vectors for word representation.

"I want a glass of Orange juice to go along with my cereal."

$x_{ij}^{st} = \# \text{ of times } j \text{ appears in the context of } i$

$$\text{minimize } \sum_{i=1}^{10k} \sum_{j=1}^{10k} f(x_{ij}) \left(\theta_i^T e_j + b_i + b_j - \log x_{ij} \right)^2$$

$\boxed{i, j = t, c}$

or $(\theta_t)^T e_c$

Correlation

$f(x_{ij}^{st})$
 this is the of at
 (stop words)

* might give more weight to
 stop words and lure on more attention. words.

if $(x_{ij}^{st} = 0)$ then $\log 0 = \text{undefined}$, so we will add
 weighing term $f(x_{ij}^{st})$, $f(x_{ij}^{st}) = 0$, if $\boxed{x_{ij}^{st} = 0}$.
 $\therefore 0 \cdot \log 0 = 0$.

f can be chosen to be heuristic to accomplish this.

Take average after calculating.

$$e_w^{(\text{final})} = \frac{e_w + \theta_w}{2}$$

APPLICATIONS USING WORD EMBEDDINGS.

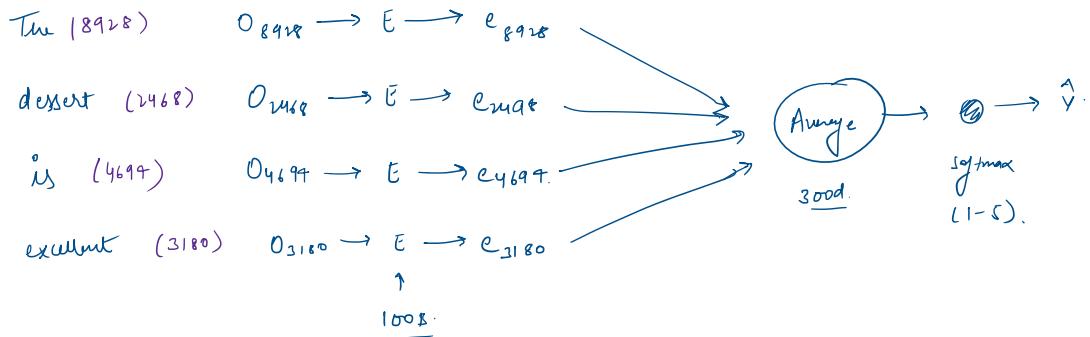
SENTIMENT CLASSIFICATION.

$$x \longrightarrow y$$

"dumb is excellent"

$s \neq$

" dessert is excellent"	5*
" Service is slow"	1*
" good for quick meal, but nothing special"	3*



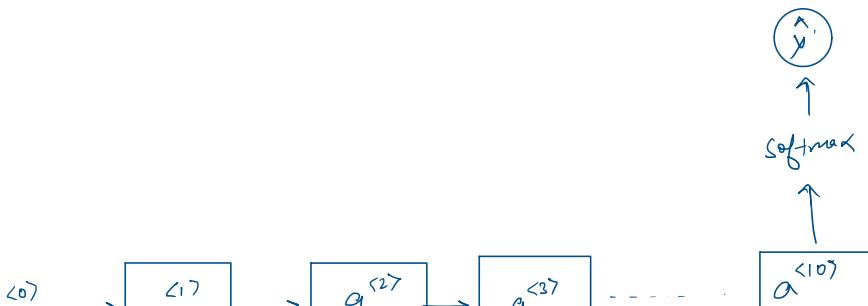
- > ignores word order.
- = " Completely lacking in good taste, good service and good ambience."
- * negative review, but "good" appears a lot!

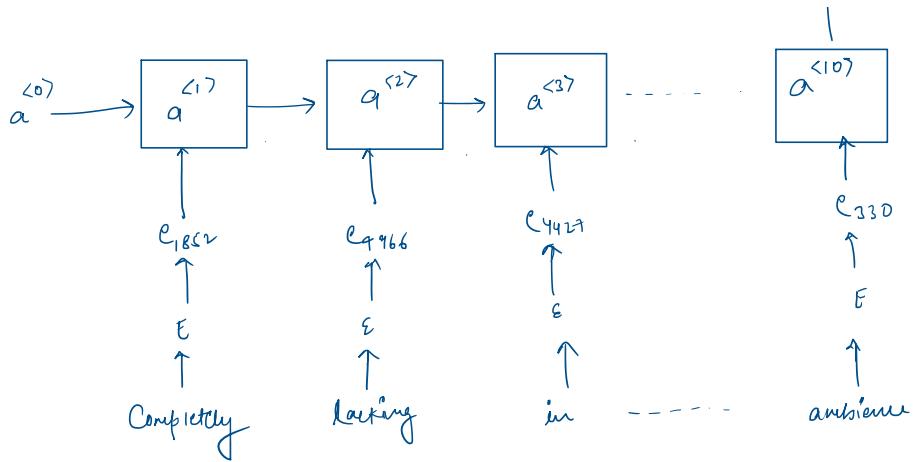
RNN for Sentiment Classification.

Completely lacking - - - - - - - - - ambience.



- ① Calculate one-hot vector.
- ② Embedding Matrix
- ③ find embedding vector. ($0 * \bar{E}$)
- ④ feed e to RNN





"Many to One?"

DEBIASING WORD EMBEDDINGS

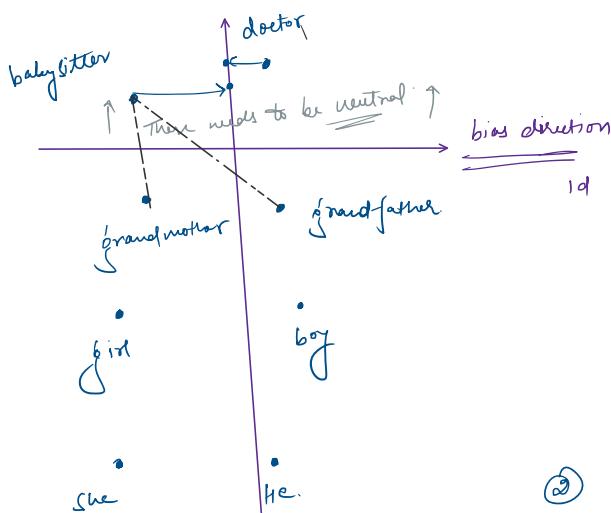
Man : Women :: King : Queen.

Man : Programmer :: Women : Homemaker. \times

Father : Doctor :: Mother : Nurse \times

* Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of text used to train the model.

non-bias direction. $\overrightarrow{w_{id}}$



① Identify bias direction

(Taking gender bias here.)

$$\begin{aligned} & e_{\text{he}} - e_{\text{she}} \\ & e_{\text{male}} - e_{\text{female}} \\ & e_{\text{boy}} - e_{\text{girl}} \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Average}$$

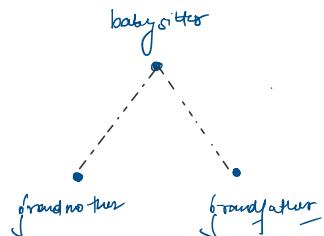
② Neutralize: For every word that is not definitional, project to get rid of bias.

③ Equalize pairs :-





* Move grandmother and grandfather at pair points equidistant from a axis



ML: VARIOUS SEQUENCE TO SEQUENCE ARCHITECTURES.

Basic Model

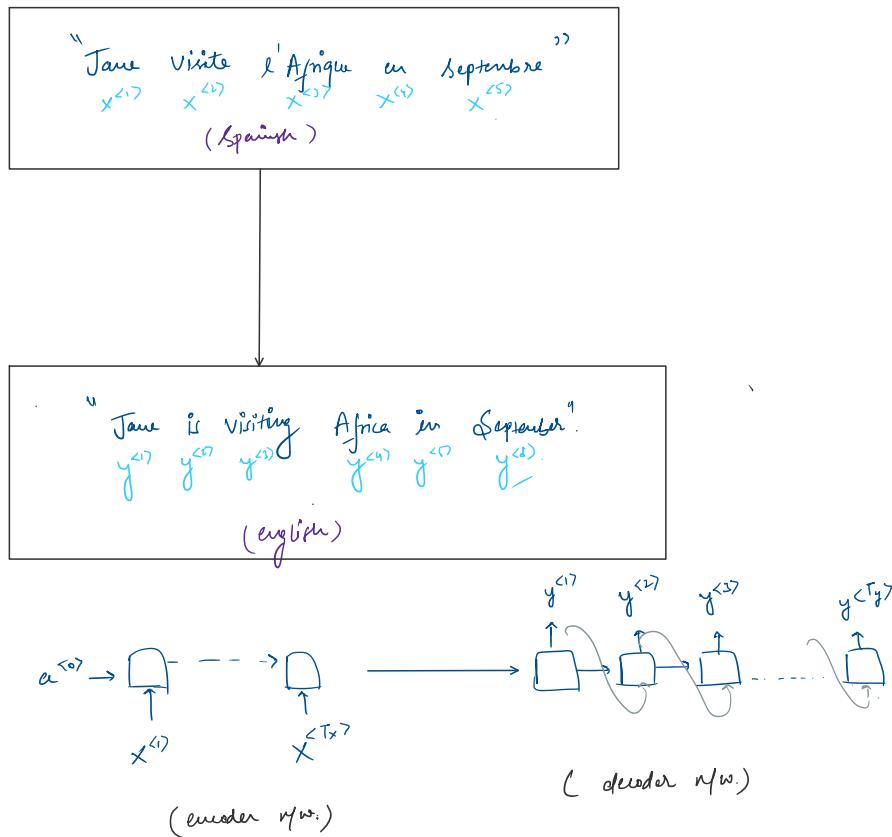
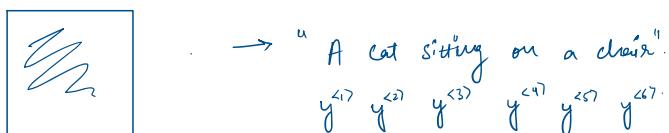
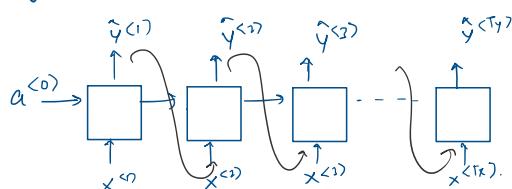


Image Captioning

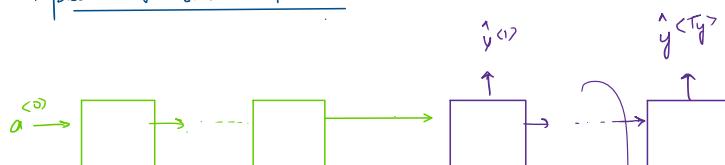


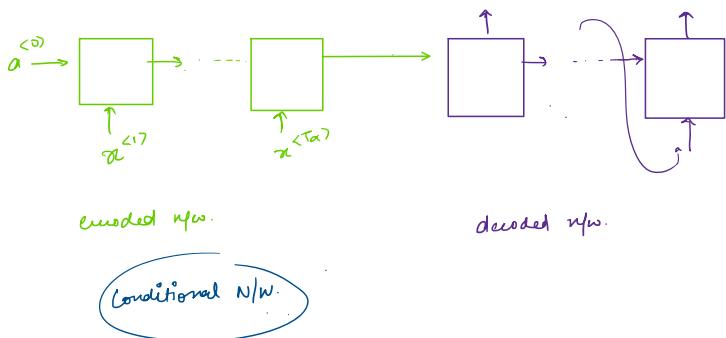
Picking the most likely sentence

Language Model (as in first week)



Machine Translation Model





"Jane visited Africa in September"

$$P(y_1, y_2, y_3, y_4, \dots)$$

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September Jane will visit Africa.
- Her African friend welcomed Jane in September.

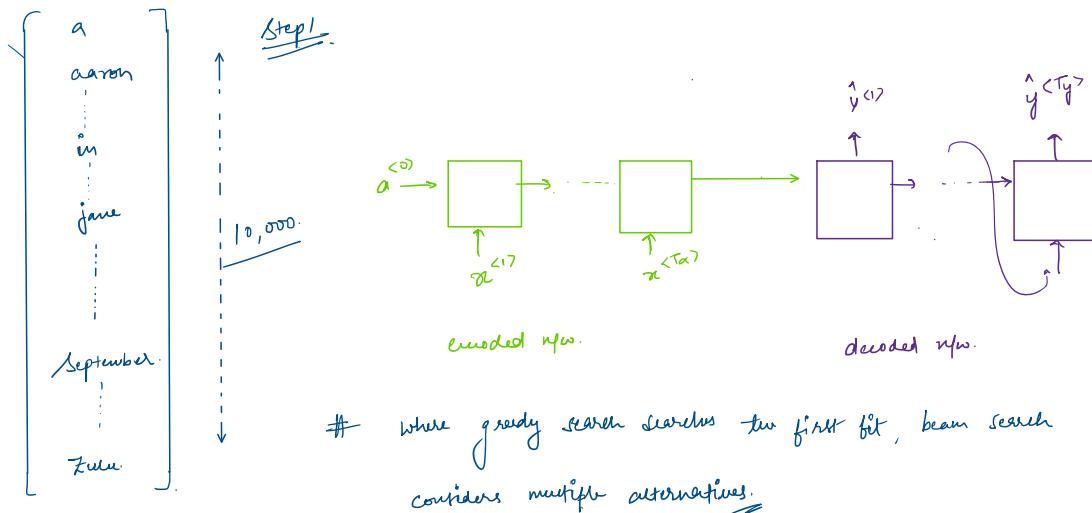
"beam search"

good translation

bad translation ==

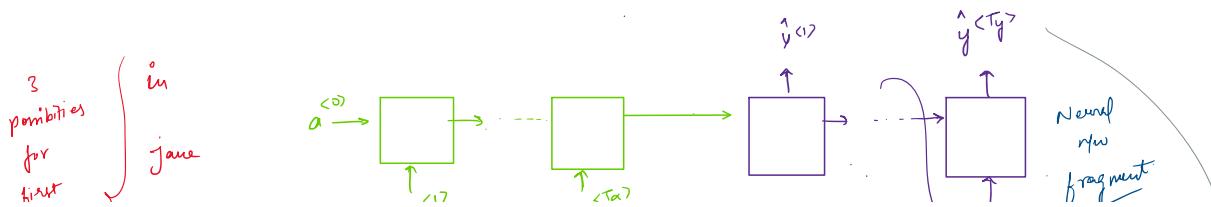
Beam Search.

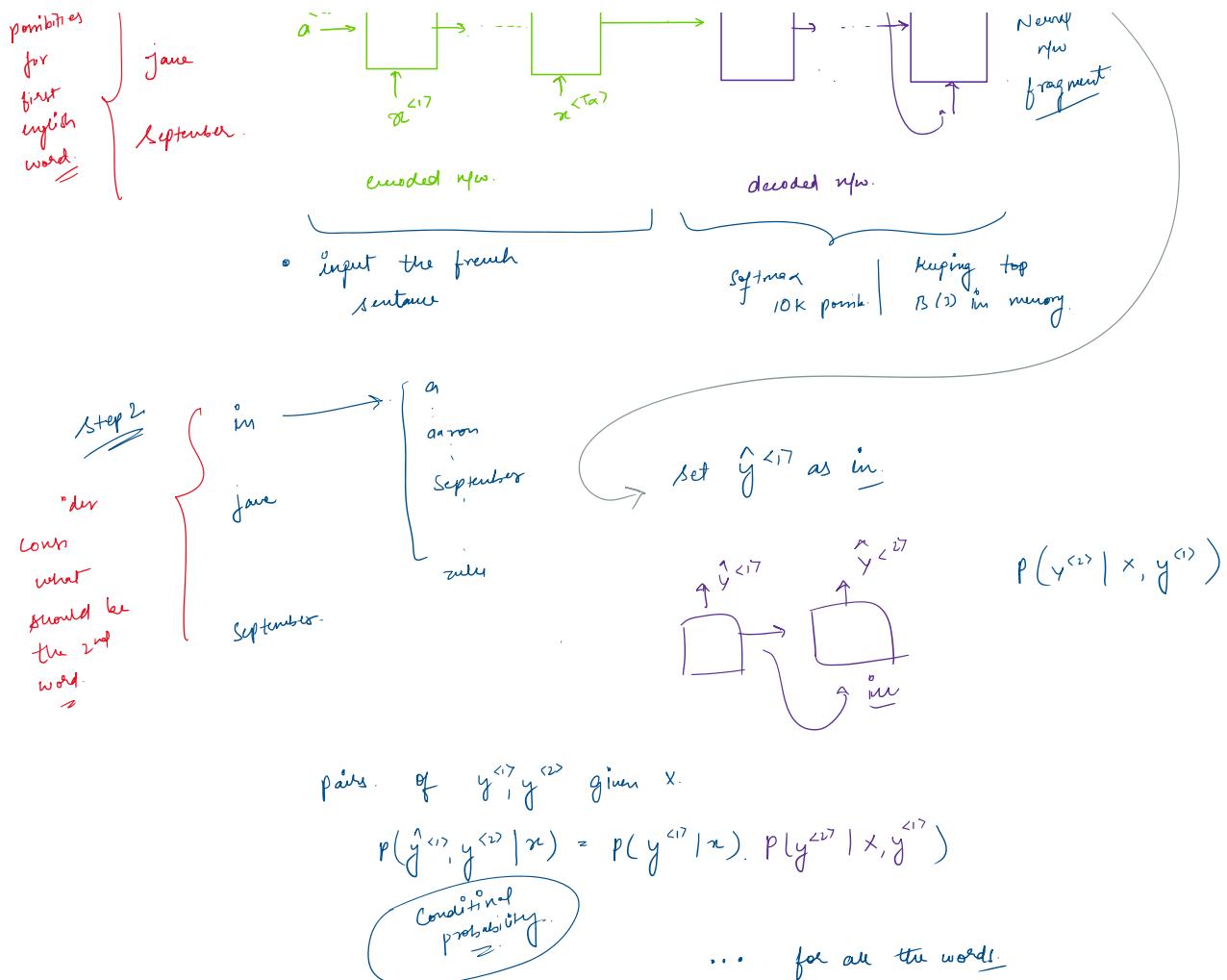
> Beam Search Algorithm



- * beam search algorithm has a perimeter B (Beam width.)

$B=3$ (say). • Considers 3 possibilities.





Refinements to Beam Search

Length Normalization

$$\text{argmax}_y \prod_{t=1}^{T_y} p(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

\log

$$p(y^{<1>} \dots y^{<T_y>} | x) = p(y^{<1>} | x) \cdot p(y^{<2>} | x, y^{<1>}) \cdot p(y^{<3>} | x, y^{<1>}, y^{<2>})$$

* Normalize by number of words.

$$\frac{1}{T_y} \sum_{t=1}^{T_y} \log p(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

Beam Width?

Larger :- Higher Accuracy, more computational.

Smaller :- Low Accuracy , Faster

Error Analysis On Beam Search.

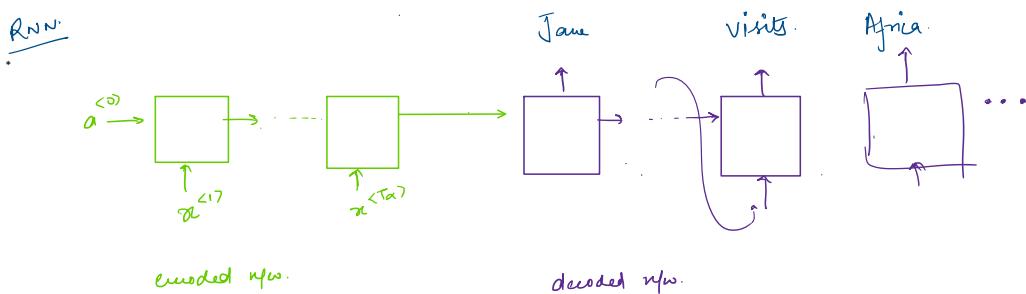
- # Beam search is an approximate search algorithm. / heuristic search algorithm.

["] Jane visite l'Afrique en septembre. ["]

Theme : Jane visits Africa in September. (y*)

Algorithm : Jane visited Africa last September. (y)

- RNN (encoder and decoder)
 - Beam Search.



$$\text{mod. } \left(p(y^* | x), \quad p(y | x) \right)$$

Jane visits Africa in September. (y*)

$$P(y^* | x)$$

Jane visited Africa last September. (iy)

$$P(\hat{y} | x)$$

Case 1 : $P(y^*) > P(\hat{y})$

\Rightarrow but beam search chooses $P(\hat{y})$.

⇒ Beam Search is at fault.

Case 2 : $P(y^*) < P(\hat{y})$

\Rightarrow y^* is a better translation than \hat{y} . But RNN predicts

$$p(y^*) < p(\hat{y})$$

⇒ RNN model is at fault.

Bleau Score

(Bilingual Evaluation Understudy Score.)

- given a french sentence there can be multiple english translation,
this can be solved with "Bleau Score"

$\alpha - \alpha - \alpha - \alpha$

French : Le chat est sur le tapis.

Refrene1 : The cat is on the mat.

Refrene2 : There is cat on the mat.

BLEAU Score measure how good a machine translation is.

$$\text{Precision} := \frac{2}{7}$$

Bleu Score on Bigrams

"pairs of words next to each other."

Refrene1 : The cat is on the mat.

Refrene2 : There is cat on the mat.

MT O/P : The cat the cat on the mat

(Machine Translation)

	MT O/P	Refrene1
the cat	Count	2
Cat the	Count	1
Cat on	Count	1
on the	Count	1
the mat.	Count	1
	CountClip	1 0 1 1 1

Modified bigram Precision :-

$$\frac{\sum \text{Count Clips}}{\sum \text{Count}}$$

$$= \frac{4}{6} = \frac{2}{3}$$

Bleu Details

P_n = Bleu Score of n -grams.

P_1, P_2, P_3, P_4

Combined bleu score = $BP \exp\left(\frac{1}{4} \sum_{n=1}^4 P_n\right)$.

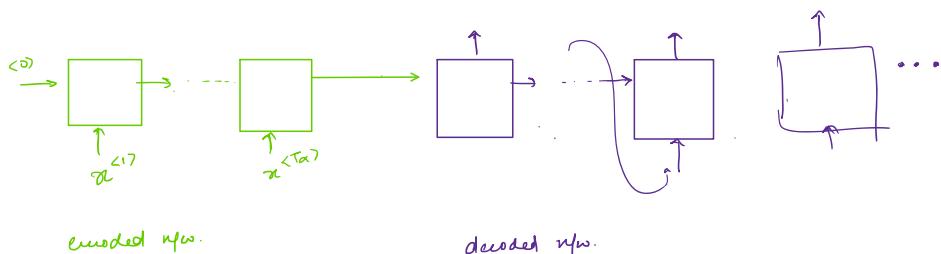
(brevity penalty)

$$BP \begin{cases} 1, & \text{if MT-Output length} > \text{reference O/P Length} \\ \exp\left(1 - \frac{\text{reference-O/P Length}}{\text{MT-O/P Length}}\right) & \end{cases}$$

Attention Model Intuition

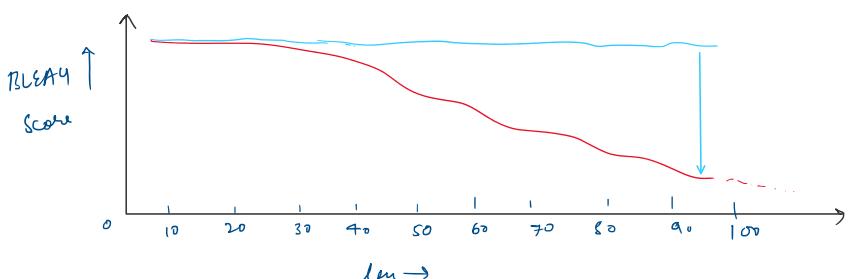
The problem of long Sequence

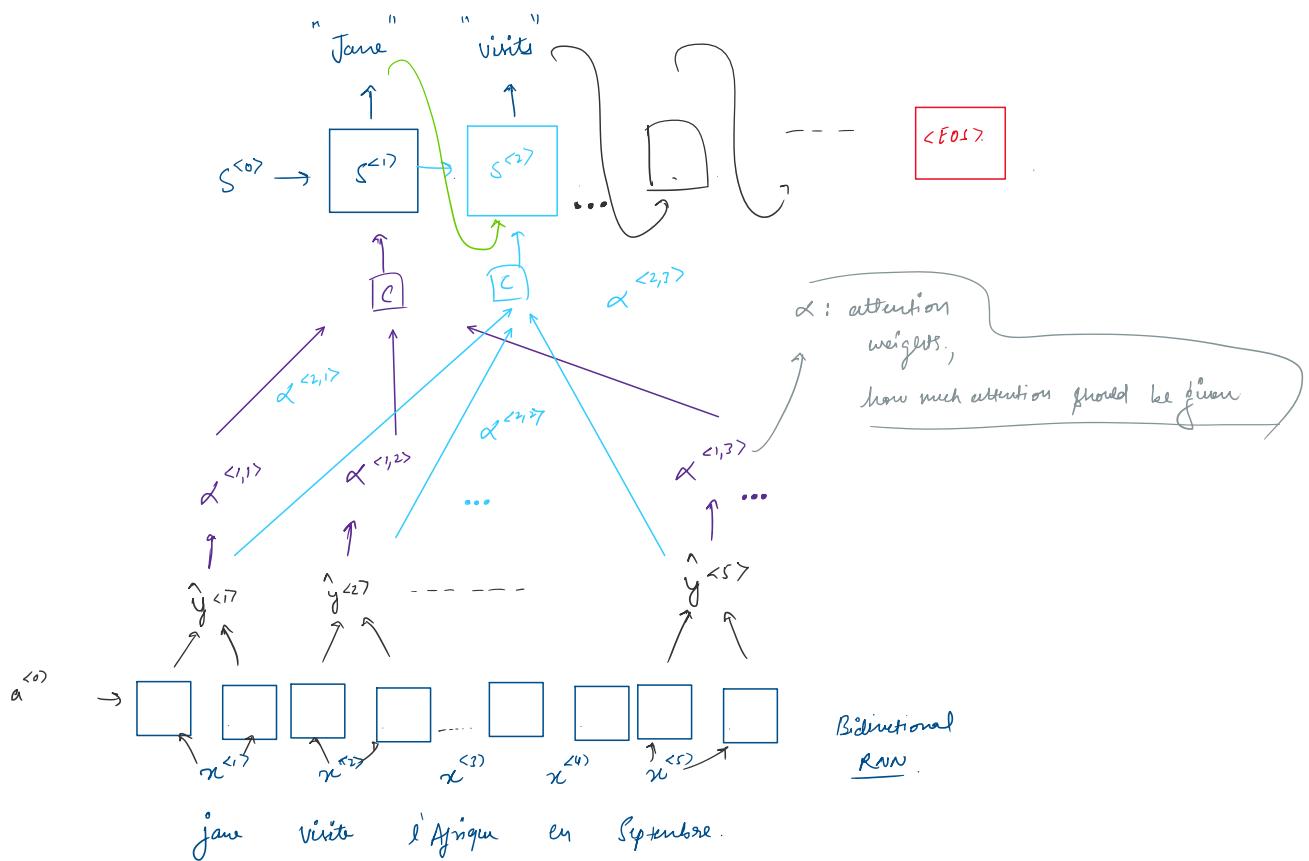
> given a long french sentence



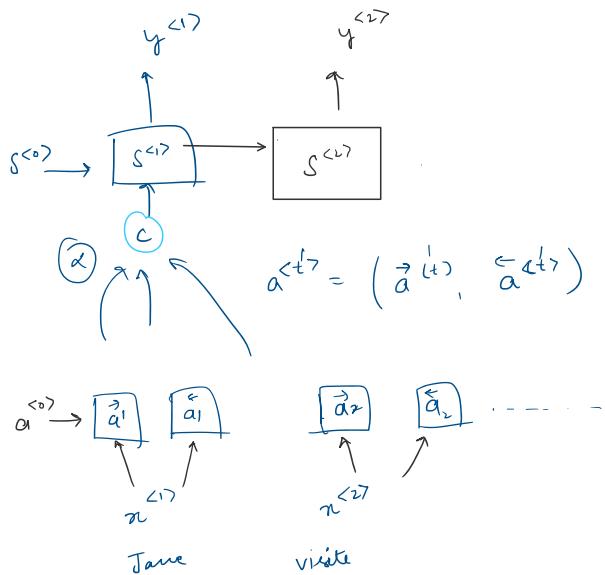
* It is hard to memorize the whole sentence, what neurons will do is

translate a part of sentence and move forward.





Attention Model



$$\sum_{x^i} \alpha^{<1,x^i>} = 1$$

$$c^{<1>} = \sum_{x^i} \alpha^{<1,x^i>} a^{<x^i>}$$

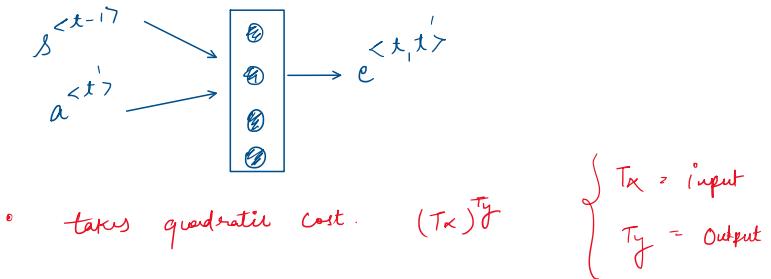
$\alpha^{<t,x^i>} = \text{amount of attention } y^{<t>} \text{ should pay to } a^{<x^i>}$

Computing Attention

$a^{<t,x^i>} = \text{amount of attention } y^{<t>} \text{ should pay to } a^{<x^i>}$

$$\alpha^{<t,x^i>} = \exp(e^{<t,x^i>})$$

$$\alpha^{(t, t')} = \frac{\exp(e^{(t, t')})}{\sum_{t'=1}^T \exp(e^{(t, t')})}$$



$$\left\{ \begin{array}{l} T_x = \text{input} \\ T_y = \text{output} \end{array} \right.$$

M2 Speech Recognition

Speech recognition

$$x \longrightarrow y$$

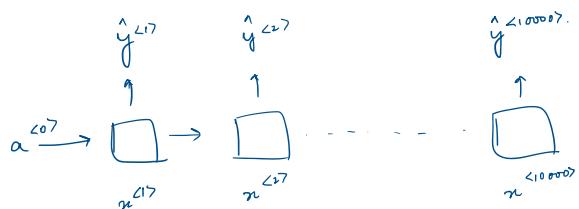
(audio clip)

(transcript)



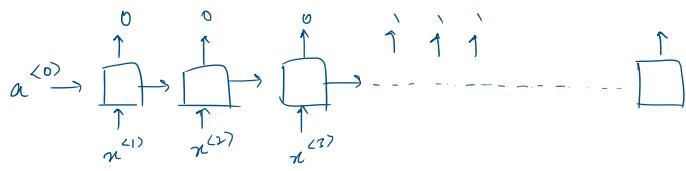
CTC Cost

"Connectionist temporal Classification"



Trigger Word detection

Spectrogram



"Hey Siri !!"

"OK google !!"

[...].